

(19) World Intellectual Property  
Organization  
International Bureau



(43) International Publication Date  
29 July 2004 (29.07.2004)

PCT

(10) International Publication Number  
**WO 2004/063769 A2**

(51) International Patent Classification<sup>7</sup>: **G01V**  
(21) International Application Number:  
PCT/US2003/041239

(22) International Filing Date:  
23 December 2003 (23.12.2003)

(25) Filing Language: English

(26) Publication Language: English

(30) Priority Data:  
60/438,259 6 January 2003 (06.01.2003) US  
10/393,641 21 March 2003 (21.03.2003) US

(71) Applicant: **HALLIBURTON ENERGY SERVICES, INC.** [US/US]; Halliburton Law Department, 3000 North Sam Houston Parkway East, Houston, TX 77032 (US).

(72) Inventors: **CHEN, Dingding**; 1705 Coit Road, #2017, Plano, TX 75075 (US). **QUIREIN, John A.**; 109 Skyline Road, Georgetown, TX 78628 (US). **WIENER, Jacky M.**; 18141 East Caley Circle, Aurora, CO 80016 (US). **GRABLE, Jeffery L.**; 6218 Laguna Del Rey, Houston, TX 77041 (US). **HAMID, Syed**; 7539 Bromwich Ct., Dallas, TX 75252 (US). **SMITH JR., Harry D.**; 12335 Kingsride, #250, Houston, TX 77024 (US).

(74) Agents: **CONLEY ROSE, PC** et al.; 5700 Granite Parkway, Suite 330, Plano, TX 75024-6615 (US).

(81) Designated States (*national*): AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BW, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NI, NO, NZ, OM, PG, PH, PL, PT, RO, RU, SC, SD, SE, SG, SK, SL, SY, TJ, TM, TN, TR, TT, TZ, UA, UG, UZ, VC, VN, YU, ZA, ZM, ZW.

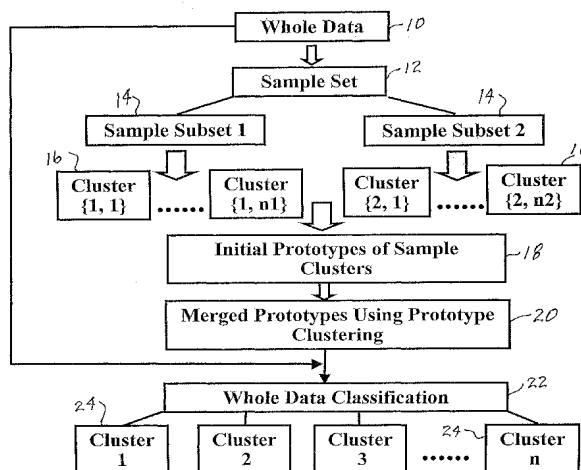
(84) Designated States (*regional*): ARIPO patent (BW, GH, GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZM, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HU, IE, IT, LU, MC, NL, PT, RO, SE, SI, SK, TR), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG).

**Declarations under Rule 4.17:**

- as to the applicant's entitlement to claim the priority of the earlier application (Rule 4.17(iii)) for all designations
- as to the applicant's entitlement to claim the priority of the earlier application (Rule 4.17(iii)) for all designations

[Continued on next page]

(54) Title: NEURAL NETWORK TRAINING DATA SELECTION USING MEMORY REDUCED CLUSTER ANALYSIS FOR FIELD MODEL DEVELOPMENT



(57) Abstract: A system and method for selecting a training data set from a set of multidimensional geophysical input data samples for training a model to predict target data. The input data may be data sets produced by a pulsed neutron logging tool at multiple depth points in a cases well. Target data may be responses of an open hole logging tool. The input data is divided into clusters. Actual target data from the training well is linked to the clusters. The linked clusters are analyzed for variance, etc. and fuzzy inference is used to select a portion of each cluster to include in a training set. The reduced set is used to train a model, such as an artificial neural network. The trained model may then be used to produce synthetic open hole logs in response to inputs of cased hole log data.

WO 2004/063769 A2



**Published:**

— without international search report and to be republished  
upon receipt of that report

*For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.*

## **NEURAL NETWORK TRAINING DATA SELECTION USING MEMORY REDUCED CLUSTER ANALYSIS FOR FIELD MODEL DEVELOPMENT**

### **CROSS-REFERENCE TO RELATED APPLICATIONS**

**[0001]** The present application claims priority from United States Provisional Patent Application 60/438,259, filed on January 6, 2003 and United States Patent Application 10/393,641, filed on March 21, 2003.

### **STATEMENT REGARDING FEDERALLY SPONSORED RESEARCH OR DEVELOPMENT**

**[0002]** Not applicable.

### **REFERENCE TO A MICROFICHE APPENDIX**

**[0003]** Not applicable.

### **FIELD OF THE INVENTION**

**[0004]** This invention relates to neural networks trained to predict a parameter in response to a plurality of inputs, and more particularly to methods of using clustering techniques and fuzzy inference to select geophysical data for use in training a neural network to produce synthetic data.

### **BACKGROUND OF THE INVENTION**

**[0005]** In the oil and gas industry today, there are several conditions that drive the need for non-traditional methods for obtaining open hole logging data. As a result, oil and gas companies are more inclined to explore such non-traditional methods for obtaining open hole logging data to help in their decision making processes. The use of cased hole logging data, in particular pulsed neutron data to generate pseudo or artificial open hole triple combo log information is one approach which has been tried.

**[0006]** One of the conditions is simple economics. Every operation carried out in a borehole takes time, which translates directly to increased cost of drilling the well. Therefore, if a logging operation in the well, e.g. an open hole log, can be avoided, it reduces the cost of drilling the well. If the same data can be obtained from another operation, e.g. a cased hole pulsed neutron log, then the actual open hole log can be skipped, saving time and money.

**[0007]** Adverse drilling conditions often make open hole logging expensive, risky or essentially impossible. Such conditions include extreme wash outs, shale bridges, caving, etc. These conditions may make it physically impossible to run an open hole logging tool

in the hole. If the tool can be run, the conditions may prevent collection of useful data in at least portions of the well.

**[0008]** Modern drilling techniques may make open hole logging risky or impossible. For example highly deviated wells may have high rates of turn or high angles which make it difficult or impossible to run an open hole tool. Some companies use slim holes, e.g. 3.5 inch diameter wells, which are too small for available open hole logging tools. However, pulsed neutron logging tools are available for running in such wells after they are cased.

**[0009]** As a result of these conditions, efforts have been made to produce synthetic or artificial open hole type logs from real data taken by pulsed neutron logging tools. However, various difficulties have been encountered in developing the predictive tools or models which are used to create such synthetic logs. For this approach to be successful, the models must produce accurate synthetic logs which can be relied on.

**[0010]** Various predictive tools have been used in processing geological logging data for many years. A field data based predictive model usually takes selected measurements of specific logging tools as inputs and produces predicted outputs using either a deterministic function or an empirical function generated from a training process. As a typical predictive framework, the artificial neural network (ANN) has received special interest and demonstrates increased use in petrophysical applications. To build an ANN model, data are selected from well logs, trained with optimization algorithms, and tested in different wells for validation. In the course of this process, data selection not only produces the greatest impact on the scope and applicability of the model, but also affects its accuracy and generalization performance. This is especially true if a single model for the field/reservoir is desired, and the data for all training wells and testing wells need to be normalized to a "field histogram". Since the uncertainty induced by different environmental factors and/or systematic errors may somehow corrupt the field data integration and pre-processing, special attention and treatment should be given to training-data selection.

**[0011]** The training-data selection is more heuristic than systematic in most neural network applications. One of the common heuristic approaches is to use a predetermined data percentage to randomly select the training, validation and testing data sets, which may cause the training results to be sensitive to the specific data splitting, especially if only single well data is available. For multiple-well training-data selection, it is quite often the case to define a resampling strategy to remove a certain amount of data in each individual well, and make the combined data set fall within a specific size limit. This procedure allows the use of some powerful, but memory-constrained training algorithms (Levenberg-Marquardt-based algorithms, for example). Otherwise, some sub-optimal training

algorithms (gradient-descent-based algorithms) must be used with sacrificed training accuracy. However, as discussed above, decision-making is difficult in determining the resampling strategy without a deep understanding of the nature of the multiple well data. Evenly scattered interval sampling (systematic sampling with respect to depth) with reduced density may remove some redundant data, but may also remove some useful information at the same time such as thin bed data.

**[0012]** There is a tendency today to integrate ANN technology with other data mining and artificial intelligence technologies for predictive model development. The advantages of using integrated technologies include enhanced predictability of the data, improved interpretability of the results, and extended applicability of the model. However, its trade-off with processing complexity should also be considered.

**[0013]** It would be desirable to have ways (1) remove faulty, redundant and insignificant data, (2) detect inconsistent data, (3) have the ability to "add", i.e., duplicate samples in key target zones.

#### SUMMARY OF THE INVENTION

**[0014]** The present invention provides improved methods for selecting training data for training a predictive model to predict target data and provides an improved trained model.

**[0015]** In one embodiment, input data is multidimensional geophysical data. The input data is divided into subsets and clustering methods are used to divide each subset into a plurality of clusters. A model or prototype is produced for each cluster. Clustering methods are applied to the prototypes to generate a second set of clusters. All original data is then assigned to the second clusters. The clusters are then analyzed to select data for use in training a model.

**[0016]** In one embodiment, target data is assigned to, or linked with, corresponding training data in each cluster. The combined data clusters are then analyzed to select data for use in training a model.

**[0017]** In one embodiment, the second clusters or combined clusters are analyzed for factors including cluster size and data variance or dispersion. Fuzzy inference is then used to select a portion of data from each cluster for inclusion in a training data set.

**[0018]** In one embodiment, the model is an artificial neural network trained to predict, or generate artificial, logging data as target data, in response to an input of actual logging data.

**[0019]** In one embodiment, the present invention provides a method of operating a field in which open hole logs are run in only a small number of wells, cased hole logs are run in all wells, a model is trained with the data collected in the logs, and the trained model is

used to generate synthetic open hole logs for the wells in which actual open hole logs were not run.

#### BRIEF DESCRIPTION OF THE DRAWINGS

**[0020]** Fig. 1 is flow chart illustrating a method of generating input data clusters according to the present invention.

**[0021]** Figs. 2A and 2B are illustrations of two input data clusters resulting from the process of Fig. 1.

**[0022]** Figs. 3A, 3B, 3C and 3D are illustrations of four data clusters including target data linked to the input data clusters.

**[0023]** Fig. 4 is a flow chart illustrating steps of cluster analysis and characterization and a fuzzy inference system for selecting portions of data from the clusters.

**[0024]** Fig. 5 is a plot of input data cluster distance mean deviations of a training well and a testing well.

**[0025]** Fig. 6 is a cross plot of cluster distance ratio versus the cluster prediction error of the testing well.

**[0026]** Figs. 7, 8 and 9 provide three comparisons of ANN computed, i.e. synthetic, open hole logs with actual triple combo measurements.

#### DETAILED DESCRIPTION OF THE INVENTION

**[0027]** For prediction of the non-linear regression problem, the method of applying cluster analysis is to locate similar data patterns in input variables, and link them with other variables to be used as desired outputs in the training process. In essence, this involves identifying patterns in the joint distribution function. This procedure allows removal of bad data and redundant data, detection of inconsistent data, and evaluation of input/output non-linearity of different clusters associated with actual geological formations. It may also provide a qualitative link between clusters and zoned facies.

**[0028]** Many clustering algorithms are available for use. See U.S. Patent 6,295,504, issued to Ye et al. on September 25, 2001, which is hereby incorporated by reference for all purposes, for an example of use of a clustering method to identify facies of geological formations based on logging measurements. To perform cluster analysis on multiple-attribute input variables, one method is to (1) find the similarity between every pair of samples by calculating distance, (2) group the samples into a binary, hierarchical tree using the distance information generated in step (1), and (3) determine where to divide the hierarchical tree into clusters according to an inconsistency setting. However, most clustering algorithms work well on small data sets containing only a few hundred samples. In multiple-well model development, the data set may contain several tens of thousands of

high-dimensional samples, making conventional approaches impossible for direct use due to an extremely large memory requirement.

**[0029]** Fig. 1 is a flow chart that shows how data clusters are generated according to the present invention. The whole data set 10 of the training well(s) may be first reduced to a sample set 12, which may, for example, contain only a certain depth in a particular well or may contain data from only a portion of the wells which have been logged. The sample set 12 is then divided into several subsets 14, and hierarchical clustering is applied to each subset respectively to reduce the memory requirement. While Fig. 1 shows two subsets 14, it is understood that the sample set 12 may be divided into more than two subsets. Since the data has been divided into smaller subsets requiring less memory, more efficient and robust clustering methods may be used. This first-level or initial grouping puts together the data patterns with predetermined inconsistency into naturally divided clusters 16.

**[0030]** For each cluster 16, a single cluster sample prototype, or mathematical representation, is calculated at step 18. Since some prototypes drawn from different subsets may be similar, a second level clustering, among the different data subset prototypes, can be performed to merge those prototypes with a lower inconsistency setting, followed by determining a new cluster prototype at step 20. Note that since the second level clustering 20 is applied to the predetermined prototypes only, the memory requirement is limited, even though the information coverage may include all of the data from all of the training well(s). Depending on the total available data size, an intermediate clustering may be needed to make the final clustering manageable.

**[0031]** The flow chart discussed in Fig. 1 is suitable for the hierarchical clustering method, in which the inconsistency coefficient is set to different values between the initial clustering and the second level clustering. The prototype of each sample cluster can be a mean vector averaged over each attribute of cluster components. It can also be an actual sample nearest to the mean vector in Euclidean distance. Other clustering methods, such as density-based methods, model-based methods and self-organizing map, may also be used depending on the nature of the data.

**[0032]** After all cluster sample prototypes are determined at step 20, each multi-attribute input vector, that is each of the multidimensional data samples in the original data set 10, is then fitted into its nearest cluster based on the distance to the cluster prototype in step 22. This results in a second set of clusters 24 which contain all of the original data 10. The cluster prototypes are then adjusted for the whole data set, that is new prototypes are generated for each of the second clusters 24.

**[0033]** Figs. 2A and 2B provide examples of two clusters with fitted samples selected from pulsed neutron inputs and prototypes. The prototypes are the heavy lines centered within the individual data set traces. These high-dimensional cluster profiles are shown in an X-Y plane, where X is the variable index: 1 for GR (gamma ray), 2 for SGIN (sigma formation intrinsic), 3 for RIN (inelastic gamma count rate ratio between detectors), 4 for RTMD (capture gamma count rate ratio between detectors), 5 for NTMD (near detector overall capture gamma count rate), 6 for FTMD (far detector overall capture gamma count rate) and 7 for SGBN (sigma borehole in near detector). Y is the normalized magnitude from -1 to +1 over the whole data range. Clearly, there is data similarity within the illustrated clusters, and data dissimilarity between the clusters.

**[0034]** The goal of performing cluster analysis, as described above, is to support artificial neural network, ANN, training-data selection. After input variables are classified into nearest clusters 24, each cluster is linked to its corresponding counterpart, the measurements to be used as targets in the training process. Most currently used plot functions cannot effectively show such comprehensive multiple-input/multiple-output relationships of the data. Cross plotting, for example, is limited by showing only a single input versus a single output. Post regression plots show correlation between the predicted outputs and the desired outputs, but the associated inputs are hard to be displayed at the same time. The conventional logging plot provides a separate curve for each individual measurement along the well depth, but does not include an ensemble of similar patterns of multiple measurements. Cluster linkage, however, conveniently provides more integrated graphic support to facilitate comprehensive analysis.

**[0035]** In Figs. 3A through 3D, four plots, Cluster A through Cluster D, are presented to link the cluster inputs of a cased-hole pulsed neutron tool with corresponding actual open-hole triple-combo measurements. The number of variables is extended to ten, with the first seven being the same input variables as shown in Figs. 2A and 2B, and last three being the triple-combo measurements (index 8 for deep resistivity, 9 for neutron porosity, and 10 for bulk density). The range of each open-hole attribute is also normalized between -1 and 1. Note that in each plot, the cluster analysis was only applied to the first seven input variables. The target data was accordingly linked with its input sample index recorded during the process. This approach provides a common framework for evaluating cluster input/output relationships, and for another important application, novelty detection, where target measurements are not available.

**[0036]** From Figs. 3A-3D, it can be seen that the input/target relationship in Cluster A is quite linear. In Cluster B, the sample dispersions of open-hole measurements are much



larger than that of the input variables, indicating a high non-linearity or inconsistency. Cluster B is a typical representative of tool responses to coal streaks that provide a problem for either pulsed neutron or density simple ratio routines. There is an outlier in Cluster C. Cluster D contains gas zone samples with small variation in both neutron porosity and bulk density. These plots provide a useful tool for problem diagnosis and model pre-assessment before the neural network is trained. These different cluster patterns may also be meaningful in developing zone or facies-based multiple models to improve prediction accuracy.

**[0037]** In the above steps, data inputs 10 have been clustered to close proximity (step 24, Fig. 1) and linked with target measurements in Figs. 3A-3D. There can be up to several hundred clusters for the training well(s) depending on reservoir characteristics and geological formation types. In the next part of the process, cluster analysis and graphic inspection are integrated to characterize each cluster. This characterization process may be followed by a decision-making system to determine how many samples need to be selected from each cluster for model development.

**[0038]** Fig. 4 shows a simplified flow chart of the cluster characterization and selection process. Each data cluster 24 is analyzed for cluster variance and other metrics at step 26 and characterized accordingly. For example, one useful metric is the cluster size, i.e. how many multidimensional data samples are included in each cluster 24, may be measured at 28. For cluster variance analysis, the distance mean, defined as the mean of the within-cluster-sample distance to the cluster prototype, is probably the most important parameter. This parameter can be expressed in the ratio form, called dispersion ratio at step 30, and used as a cluster non-linearity index. The dispersion ratio is the ratio of target distance mean divided by the input distance mean. The zone indicator at step 32 takes account of the weight of key zones, allowing duplication of the key samples in the training set. All these parameters, plus others if necessary, may be used as inputs of the decision-making system for training-data selection.

**[0039]** Fuzzy inference is the process of formulating the mapping from a given input to an output using fuzzy logic, and has found many applications in decision-making. In a preferred embodiment, a fuzzy inference system is used to receive inputs from the cluster characterization process, and produce an output equal to the percentage of data to be used for training in each cluster. The inputs may be relative cluster size 28, cluster dispersion ratio 30, key zone indicator 32 and/or other variables derived from the preceding process. The functionality of the fuzzy inference system 34 can be described in several steps. In the step of fuzzification, the system receives the crisp input of each

variable, and converts it to fuzzy input, which is a degree of satisfaction defined by the adaptive membership function. In the step of logic-operation, logical AND and OR operations are performed to represent the antecedent of each fuzzy rule, and its consequence is obtained in the step of implication. The outputs of the multiple rules are then aggregated in the step of aggregation to form a fuzzy set. Finally, for a given cluster, a single number is calculated at step 36 from the defuzzification step to indicate a percentage of the samples to be used in the training set.

**[0040]** In this embodiment, fuzzy rules constitute the basic decision-making strategy. The implementation of the other steps is straightforward. For example, the character A may represent relative cluster size, B may represent cluster dispersion ratio, C may represent key zone index ranged from 0 to 1, and D may represent percent data to be selected. An example of fuzzy rules may be:

If A is *large*, B is *low* and C is *low*, then D is *low*;

If A is *small*, B is *med* and C is *high*, then D is *high* ;

etc.

where the membership function *large*, *small*, *low*, *med*, *high* should be defined for each variable involved. To make the fuzzy inference system a useful tool, an adequate number of fuzzy rules are required. Similar to data preparation for other predictive tools, flexibility always exists for ANN training-data selection due to the capacity of the computer, experience of the designer, and limited information source. Some heuristic approaches can be combined with cluster information to obtain the best compromise between the processing simplicity and model predictability.

**[0041]** The percentages determined at step 36 are used to select a portion of data from each cluster 24. The data selected from each cluster 24 is combined into a training data set to be used to train a model for predicting the target data from real input data. As indicated above, a preferred embodiment uses an artificial neural network, ANN, as the predictive model. Normal training methods are used to train the ANN. For example, the training data set may typically be split into training and validation subsets. However, the training data selection process of the present invention results in a reduced training set which allows use of preferred training algorithms for training an ANN. After training and validation, the ANN may be tested with input data from other wells to determine if the model is good enough. When a model has been shown to accurately predict, or generate artificial logs, e.g. open hole logs, from real input data, for example cased hole log data, then it may be used to generate such artificial logs for other wells

**[0042]** In the process described above, it was assumed that multiple well data were available, and all the available data was analyzed and processed at the same time. It is quite often the case that the original ANN model is trained with certain well(s) data first, and tested on a different well later to determine if the model is good enough. If not, the designer may want to add some new data selected from the testing well to the previous training set and retrain the network without re-processing the whole data set of the multiple wells. An example of how the above described process can be used to simplify retraining with the additional data is described below.

**[0043]** Before deciding whether or not to add the data from another well to the training set, novelty testing should be performed first, using the previously generated cluster prototypes from the training well(s) to classify the new data from the testing well. The analysis characterizes the new data and indicates how the testing inputs are similar to the training inputs, and how this similarity is related to the prediction error of testing data. The analysis can be used to establish a criterion to help add only “novelty” for new training, and can improve the applicability of the field model. This reduces the total amount of data added to the training set so that the most effective training algorithms can still be used.

**[0044]** Figs. 5 and 6 summarize some results of this example. As previously discussed, the ANN model takes seven measurements from the cased-hole pulsed neutron logs as inputs to predict open-hole triple-combo outputs. Fig. 5 is a plot of input cluster distance means of a training well, the lower curve, and a testing well, the upper curve. The clusters, 189 in total number for the training well, were statistically generated using about 4500 feet of log data. The testing well, which is about three miles away from the training well, has its data (about 5000 feet of log data) fitted to 182 of the training well clusters, leaving blanks in the plot for the non-filled clusters. It can be seen from Fig. 5 that for each cluster, the distance mean of the testing well is consistently larger than that of the training well. Highly deviated mean values between the wells often indicate the existence of over range cluster inputs in the testing well.

**[0045]** Fig. 6 depicts cluster distance ratio, which is the ratio of the input cluster distance of the testing well over the same cluster distance of the training well, as an input similarity index, and plots that ratio versus the cluster prediction error of the testing well. The measurements in Fig. 6 are divided into four quadrants along the prediction error (root of mean-squared-error) set point of 0.2 and the distance ratio set point of 2, assuming that those values can be used as simple thresholds. In quadrant I, out-of-boundary data show that significantly different inputs lead to larger prediction error. In contrast to quadrant I, 138 of 182 clusters in quadrant III, which is the dominant part of the whole data set, are

statistically located in the region where the clustered inputs are similar to their training well counterparts, and the prediction error is low. Probably the most questionable quadrant is quadrant II, which contains clusters with higher uncertainty exhibited (higher prediction error). The reasons for the coexistence of lower input dissimilarity and higher prediction error could be that dissimilar inputs are classified into the same cluster due to the use of a single distance measurement, or, non-linearity involved is so high that variation observed in output cannot be differentiated by its input. Any bad data in cased-hole and/or open-hole logs, and any inappropriate pre-processing of inputs and outputs will also affect prediction accuracy. As contrasted with quadrant II, quadrant IV contains the clusters whose cased-hole pulsed neutron inputs are linearly correlated with corresponding open-hole triple-combo outputs. The prediction is therefore adequately accurate even though the testing inputs seem out of range compared with the training inputs.

**[0046]** Based on the analysis stated in the previous paragraph, the clusters in quadrant I are clearly novel, and should be selected as retraining candidates. The major part of the data in quadrant III can be excluded from retraining consideration in general because of its higher input similarity and lower prediction error. Data in quadrant II need to be reinvestigated with caution to determine what causes inconsistency. Graphical inspection of cluster-analysis results can help identify problems in this quadrant. Data in quadrant IV is not crucial due to its linearity with output. Finally, only about 20 percent of the data in the second well (mainly from 40 clusters) were combined with the data of the primary training well to build the multi-well model.

**[0047]** Figs. 7, 8 and 9 provide three examples of ANN computed, i.e. synthetic, open hole logs plotted with actual triple combo measurements. The neural networks used in these examples were constructed by two layers (one-hidden layer), seven inputs, and three outputs.

**[0048]** Fig. 7 displays ANN predictions of 250-ft open-hole log data (including several gas bearing intervals, such as 1105-40 ft) versus the actual log data of a first well, using the model trained from the same single-well data. About 50 percent of data in the first well were used in the training. The post-regression coefficients between the ANN predictions and the actual measurements over the 4500-ft logs can be up to 0.86 for deep resistivity, 0.96 for neutron porosity and 0.95 for formation density. Excellent agreement is observed between the actual open-hole logs and those computed from the pulsed neutron data.

**[0049]** In Fig. 8, the same model was tested on a second well. Data from this second well was not used to train the ANN previously developed. This well also contained gas

zones, including the interval 1120-70 ft. Agreement between the logs is good, but not as good as in Fig. 7.

**[0050]** The testing in Fig. 9 is on the same log subset of the second well as shown in Fig. 8, but the model was trained with significantly reduced joint data set of two wells using methods described above. Using only about 25 percent of data from each well, the overall (including bad data) post-regression correlation coefficients of the multi-well model on the second well were improved from 0.67 to 0.78 for deep resistivity, from 0.88 to 0.92 for neutron porosity, and from 0.85 to 0.91 for bulk density. However, the prediction accuracy of the first well was sacrificed slightly to balance total error due to the data inconsistency between the two wells. It is apparent from Fig. 9 that it is possible in a multi-well environment to use neural nets and clustering concepts to accurately simulate open hole triple combo logs from pulsed neutron log data. Also note that in Figs. 7, 8 and 9 the gas zones were accurately profiled on the computed density logs. This is especially noteworthy since pulsed neutron tools do not contain gamma ray sources (such as those present in all density logging tools), and hence the gas zones represented a very challenging environment for the ANN model.

**[0051]** The above description and drawings illustrate how cluster analysis can be integrated with graphical visualization and fuzzy decision making to support sample selection in field model development using a neural network as a predictive framework. The methods discussed can also be used to support other analyst-based data interpretation and problem diagnosis with different predictive tools, i.e. other models. This approach greatly improves transparency of the conventional "black box" neural network to the designer, extends the model utility from the single-well source to multi-well sources in a cost-effective manner, and provides a powerful means to evaluate the data processing, input/output selection and the tool limitation for goal-related data mining. In logging and petrophysics applications, this method is most suitable to support multi-well field model development for medium-to-large-sized high-dimensional data interpretation. Using this method it can be seen that it is possible in a multi-well environment to generate excellent open-hole triple combo logs from cased-hole pulsed neutron data.

**[0052]** The following steps outline an embodiment of the cluster-analysis-based fuzzy reference system for neural network training sample selection according to the present invention.

**[0053]** 1. Normalize the pre-processed multi-dimensional data.

**[0054]** 2. Partition the well data into several subsets.

**[0055]** 3. Evenly sample each subset along the coordinate of well depth.

- [0056]** 4. Find the input sample cluster with predetermined inconsistency coefficient.
- [0057]** 5. Locate the prototype of each sample cluster.
- [0058]** 6. Merge sample prototype by recluster with lower inconsistency coefficient.
- [0059]** 7. Relocate the prototype of each sample cluster.
- [0060]** 8. Fit all data into its nearest sample cluster.
- [0061]** 9. Link each cluster's input and target data to be used in regression model.
- [0062]** 10. Perform cluster statistics.
- [0063]** 11. Determine the range of membership function of the fuzzy inference system.
- [0064]** 12. Characterize the fuzzy input of each cluster.
- [0065]** 13. Calculate fuzzy system output to obtain data percentage to be chosen from each cluster.
- [0066]** 14. Sample each cluster to form neural network training set.

#### ALTERNATIVE APPLICATIONS

**[0067]** The present invention has been described primarily with respect to using multidimensional pulsed neutron log data from cased wells to predict geological values normally measured by logging open boreholes. However, it has other applications. It is generally applicable to training and use of predictive models having multiple geological and/or geophysical data inputs and producing one or more geological and/or geophysical values as output(s).

**[0068]** The present invention is useful in detecting changes in the formations which occur over time due to production of oil and gas. The interfaces between water, oil and gas changes as these materials are produced. In wells which were open hole logged before production, open hole logs would be different if they could be taken after production. The present invention allows synthetic open hole logs to be generated from cased hole logs taken after production so that a comparison can be made to determine changes which result from production.

**[0069]** The pulsed neutron logging tool used in the preferred embodiments provides at least seven separate data outputs. In a large field it may be desirable to run this logging tool in only some of the wells and use a simpler and less expensive tool in the remaining wells. The simpler tool may measure some, but not all of the parameters measured by the larger tool. The full set of measurements taken in a few wells may be broken into input and target values. The input values would be only the values which the simpler tool will measure in the remaining wells in the field. A predictive model can be trained as described

in the present invention to generate synthetic logs of the target values for the wells in which only the simpler logging tool is run.

**[0070]** A similar application provides reconstruction of open hole or cased hole logs which have missing or defective data. For example, due to poor well conditions, open hole logs may have certain depth intervals without data or with defective data. After cased hole data is collected in such wells, a process like that described with reference to Fig. 7 can be used to reconstruct the open hole log or fill in the missing or defective data. This can be done by using the good open hole data, i.e. from depth intervals other than those with no data or defective data, together with cased hole data from the same intervals to train a model in accordance with the present invention. Then cased hole data from the intervals with no data or defective data can be input to the model to produce the missing open hole data.

**[0071]** In other cases, open hole data, possibly combined with cased hole data, can be used to reconstruct open hole logs with missing or defective data. Open hole logs normally produce a plurality of parameters. In some cases, poor well conditions may affect only one or two of the parameter readings. The log may include good data for the other parameters. In such a case, the parameters with good data can be used as the inputs for model training, and the parameter(s) which are partially missing or defective can be used as target data for training. The good parameters in the zones with defective data can then be input to the trained model to provide synthetic values for the missing parameters in the defective zones. In this scenario, parameters measured by cased hole logging may also be included as inputs during training and during data reconstruction if desired.

**[0072]** The data reconstruction process can also be used to reconstruct or fill in missing or defective cased hole log data. The process can be like any of those described in the preceding two paragraphs. The difference would be that the cased hole log may be missing data for one or more parameters in some depth intervals. The parameters which were accurately measured in those intervals, possibly combined with open hole parameters in those intervals, if available, can be used as inputs for training a model. The good portions of the parameter(s) which are partially missing would be the target data for training. The good cased hole data, and corresponding open hole data, if available, can then be input to the trained model to produce the missing data.

**[0073]** In some cases, open hole logs may have been run in a large number of wells in a field. At a later time it may be desired to run more current logging tools in the wells which are now cased. The process of the preferred embodiment may be essentially reversed to use the open hole measurements to predict some or all of the desired new

logging measurements, e.g. the suite of pulsed neutron log data. This can be done by running a pulsed neutron tool in some of the wells in the field and using the data collected as the target data for training a predictive model. The input data would be the original open hole logs from these same wells. The open hole logs from the remaining wells may then be used with the model to predict the pulsed neutron tool data without actually running the tool in all of the remaining wells. Alternatively, a simpler pulsed neutron tool, as discussed above may be run in the remaining wells and its measurements may be used with the original open hole measurements as inputs to a model to predict the remaining data.

**[0074]** The multidimensional input data need not be a suite of measurements taken by a single instrument or set of instruments run in a borehole at the same time. Measurements from two or more instruments in the same well may be depth correlated and combined to form a set of input data parameters. These measurements may include, among others:

**[0075]** Nuclear Magnetic Resonance (NMR)

**[0076]** Dipole Sonic

**[0077]** Electric Micro Imaging Log

**[0078]** Pulsed Neutron

**[0079]** Pulsed Neutron & Carbon Oxygen

**[0080]** Open Hole Logs

**[0081]** Open Hole Triple Combo (Resistivity, Density, Neutron, Sonic)

**[0082]** Cased Hole Production Logs

**[0083]** Subsurface Core Data

**[0084]** Formation Pressure Data

**[0085]** Vertical Seismic Profiling

**[0086]** Other types of data such as measurements of formation samples, e.g. drill cuttings or sidewall cores, may also be used. The input data may also include measurements taken from the earth's surface, e.g. seismic data, which may be depth correlated with borehole log data. In similar fashion, the actual target data used in training may be data from more than one logging tool and/or may include non-borehole data such as seismic data. Predicted target measurements may likewise be the types of measurements normally measured by a logging tool or data normally measured by other means, e.g. seismic.

**[0087]** Another alternative application for clustering techniques (including Self Organization Mapping) as applied to subsurface and/or surface measurements is facies



identification. Facies identification from logging data can be an extremely important predictive product from these techniques in terms of rock typing for fracture stimulation design, petrophysical analysis, permeability determination for fluid flow characteristics, and understanding subsurface reservoir properties. Clustering of multidimensional input data measurements will group/organize these measurements in such a way to have geologic significance and thus further the operator's knowledge of their producing reservoirs.

**[0088]** The present invention provides new flexibility in development of hydrocarbon, e.g. oil and gas, bearing fields. For example, a plan for development of a field may call for drilling a plurality, e.g. fifty, wells into the producing formations in the field. For various reasons, such as those discussed in the background section above, it may be very desirable to limit open hole logging to only a portion of the fifty wells, e.g. maybe only ten or fewer wells. After the wells are drilled and cased, cased hole logs, e.g. pulsed neutron logs, may be run in all wells in the field. The few open hole logs together with the cased hole logs from the same wells may then be used as the training data according to the present invention to develop a model, e.g. an artificial neural network, which is representative of the entire field. The model may then be used to produce synthetic open hole logs for all wells in the field, or at least those which did not have actual open hole logs, by inputting the cased hole log data into the model.

#### ADVANTAGES OF THE INVENTION

**[0089]** The present invention applies a goal related clustering method. In this new approach, the objective was not to develop a theoretically novel clustering method. Instead, it was to select and integrate clustering methods to achieve a particular goal. This invention provides a cluster-analysis-based algorithm to efficiently locate the similar data patterns, and produces results that are interpretable, comprehensible, and usable for neural network training sample selection.

**[0090]** The present invention provides a reduced memory requirement. To make the clustering algorithm suitable to large data sets, data partitioning and prototype merging methods are included in this invention and results in reduced memory requirement. The strategy of NN training sample selection can then be deliberately determined from the cluster density distribution and importance of the data patterns. This facilitates the removal of redundant data and insignificant data, and allows applying some powerful, but memory-constrained training algorithms to the well-selected data set for field model development.

**[0091]** The invention provides enhanced diagnostic capability. Problem diagnosis is challenging in data mining. In this invention, the high-dimensional data is displayed in an X-Y plane to show what the input/output mapping relationship in the cluster looks like. This

makes model pre-assessment convenient before the neural network is trained, facilitates the user to locate the outlier, to reselect input parameters and to reinvestigate the pre-processing method. In addition, it makes novelty detection on new data practical by using the same cluster analysis framework.

**[0092]** The present invention incorporates fuzzy-adapted decision making. Unlike the common practice of placing the fixed percent of total data in the training set, the new approach involves a fuzzy inference system to help decision making in training sample selection. The membership function of the fuzzy inputs and outputs are defined based on the statistical results of cluster analysis, which is problem dependent and can be adapted dynamically when the cluster statistics are changed.

**[0093]** It is apparent that various changes can be made in the apparatus and methods disclosed herein, without departing from the scope of the invention as defined by the appended claims.

## CLAIMS:

What we claim is:

1. A method for producing a training data set from a set of multidimensional geophysical input data samples for training a model to predict target data, comprising:
  - dividing a set of geophysical input data samples into a plurality of first subsets of input data samples,
  - dividing each of the first subsets into a plurality of first clusters,
  - generating a first set of prototypes each representing one of the first clusters,
  - and
  - dividing the first set of prototypes into a plurality of second clusters.
2. A method according to Claim 1, further comprising generating a second set of prototypes each representing one of the second clusters.
3. A method according to Claim 1, further comprising assigning each of the geophysical input data samples to one of the second clusters.
4. A method according to Claim 3, further comprising linking actual target data to each of the geophysical data samples.
5. A method according to Claim 4, further comprising performing a statistical analysis of data samples and target data within each cluster.
6. A method according to Claim 5, further comprising processing the results of the statistical analysis with fuzzy inference to assign a percentage to each of the second clusters.
7. A method according to Claim 6, further comprising selecting the assigned percentage of data samples and target data from each cluster to form a training data set.
8. A method according to Claim 4, further comprising generating a plot of the data samples and target data for each cluster.
9. A method according to Claim 8, further comprising visually inspecting each plot.
10. A method according to Claim 9, further comprising selecting data samples and target data to be included in or excluded from the training set.

11. A method according to Claim 5, wherein the statistical analysis includes determining one or more of cluster size, data dispersion ratio, and zone to which the data samples relate.
12. A method according to Claim 1, wherein the geophysical data comprises output readings from a pulsed neutron logging tool in at least one cased borehole.
13. A method according to Claim 1, wherein the target data comprises logging measurements taken in an open borehole.
14. A method according to Claim 13, wherein the target data comprises measurements representing one or more of neutron porosity, formation density and deep resistivity.
15. A method according to Claim 1, wherein the model is an artificial neural network adapted to predict target data in response to geophysical input data samples.
16. A method for producing a training data set from a set of multidimensional geophysical input data samples for training a model to predict target data, comprising:
  - dividing multidimensional geophysical input data samples into a set of clusters,
  - linking each multidimensional geophysical input data sample with corresponding target data, and
  - performing an analysis of the input samples and target data in each cluster.
17. A method according to Claim 16, further comprising selecting a percentage of data samples and corresponding target data from each cluster based on results of the analysis.
18. A method according to Claim 17, further comprising combining the data selected from each cluster to form a training data set.
19. A method according to Claim 16, further comprising generating a plot of the data samples and target data for each cluster.
20. A method according to Claim 19, wherein the step of performing an analysis comprises visually inspecting each plot.
21. A method according to Claim 16, wherein the step of performing an analysis comprises calculating the dispersion ratio of the data samples and target data for each cluster.

22. A method according to Claim 21, further comprising processing the results of the analysis with fuzzy inference to assign a percentage to each of the second clusters.
23. A method according to Claim 16, wherein the input data comprises output readings from a pulsed neutron logging tool.
24. A method according to Claim 23, wherein the output readings are taken in a cased borehole.
25. A method according to Claim 16, wherein the target data comprises logging measurements taken in an open borehole.
26. A method according to Claim 16, wherein the target data comprises measurements representing one or more of neutron porosity, formation density and deep resistivity.
27. A method for predicting open borehole logging measurements from actual cased borehole logging measurements, comprising:
  - collecting open hole logging measurements in a borehole,
  - collecting cased borehole logging measurements in the borehole,
  - dividing the cased borehole logging measurements into a set of clusters,
  - linking each cased borehole logging measurement with corresponding open hole logging measurements,
  - performing an analysis of the cased borehole logging measurements and corresponding open hole logging measurements for each cluster,
  - selecting a percentage of the cased borehole logging measurements and corresponding open hole logging measurements from each cluster based on results of the analyses,
  - training a predictive model with the selected measurements, and
  - using the trained predictive model to predict open hole logging measurements in response to cased borehole logging measurements.
28. A method according to Claim 27, wherein the step of performing an analysis comprises:
  - plotting the cased borehole logging measurements and corresponding open hole logging measurements for each cluster,
  - visually inspecting each plot, and
  - selecting data from each cluster based on the visual inspection.

29. A method according to Claim 27, wherein the step of performing an analysis comprises performing a statistical analysis of the cased borehole logging measurements and corresponding open hole logging measurements within each cluster.
30. A method according to Claim 29, further comprising processing the results of the statistical analysis with fuzzy inference to assign a percentage to each of the second clusters.
31. A method according to Claim 27, wherein the step of dividing the cased borehole logging measurements into a set of clusters comprises:
  - dividing the cased borehole logging measurements into a plurality of first subsets,
  - dividing each of the first subsets into a plurality of first clusters,
  - generating a first set of prototypes each representing one of the first clusters,
  - and
  - dividing the first set of prototypes into a plurality of second clusters.
32. A method according to Claim 31, further comprising assigning each of the cased borehole logging measurements to one of the second clusters.
33. A method according to Claim 27, wherein the predictive model is an artificial neural network.
34. A method according to Claim 31, wherein the cased borehole logging measurements are outputs of a pulsed neutron logging tool.
35. A method according to Claim 27, wherein the open borehole logging measurements comprise measurements representing one or more of neutron porosity, formation density and deep resistivity.

36. A method for predicting open borehole geophysical measurements from actual cased borehole geophysical measurements, comprising:
- collecting open hole geophysical measurements in a borehole,
  - collecting cased borehole geophysical measurements in the borehole,
  - selecting a percentage of the cased borehole measurements and corresponding open hole measurements as a training data set,
  - training a predictive model with the selected measurements, and
  - using the trained predictive model to predict open hole geophysical measurements in response to cased borehole geophysical measurements.
37. A method for predicting cased borehole geophysical measurements from actual open borehole geophysical measurements, comprising:
- collecting open hole geophysical measurements in a borehole,
  - collecting cased borehole geophysical measurements in the borehole,
  - selecting a percentage of the open hole measurements and corresponding cased borehole measurements as a training data set,
  - training a predictive model with the selected measurements, and
  - using the trained predictive model to predict cased hole geophysical measurements in response to open borehole geophysical measurements.
38. A method for producing a synthetic log of at least one geophysical parameter for a well, comprising:
- collecting a first log of a plurality of geophysical parameters, including the at least one geophysical parameter, in a first well, the log comprising a plurality of multidimensional data samples,
  - dividing the data samples into a set of clusters based on the geophysical parameters other than the at least one geophysical parameter,
  - selecting data from each cluster,
  - training a predictive model with the selected data,
  - collecting a second log of the plurality of geophysical parameters, excluding the at least one geophysical parameter, in a second well, and
  - inputting the second log to the predictive model to produce a synthetic log of the at least one geophysical parameter for the second well.

39. A method according to Claim 38, further comprising analyzing the data in each cluster.
40. A method according to Claim 39, further comprising performing an analysis on data in each cluster.
41. A method according to Claim 39, further comprising:
  - plotting the data in each cluster, and
  - visually inspecting the data plots.
42. A method according to Claim 41, further comprising identifying formation type represented by a cluster.
43. A method for producing a synthetic value of at least one geophysical parameter for a well, comprising:
  - collecting a first data sample set of a plurality of geophysical parameters, including the at least one geophysical parameter, relating to a first well,
  - dividing the first data sample set into a set of clusters based on the geophysical parameters other than the at least one geophysical parameter,
  - selecting data from each cluster,
  - training a predictive model with the selected data,
  - collecting a second data sample set of the plurality of geophysical parameters, excluding the at least one geophysical parameter, relating to a second well, and
  - inputting the second data sample set to the predictive model to produce a synthetic value of the at least one geophysical parameter for the second well.
44. A method according to Claim 43, wherein the first data sample set comprises geophysical parameters normally measured by open hole logging and by cased hole logging of the well.
45. A method according to Claim 43, wherein the first data sample set comprises seismic data.
46. A method according to Claim 43, wherein the first data sample set comprises sidewall core data.

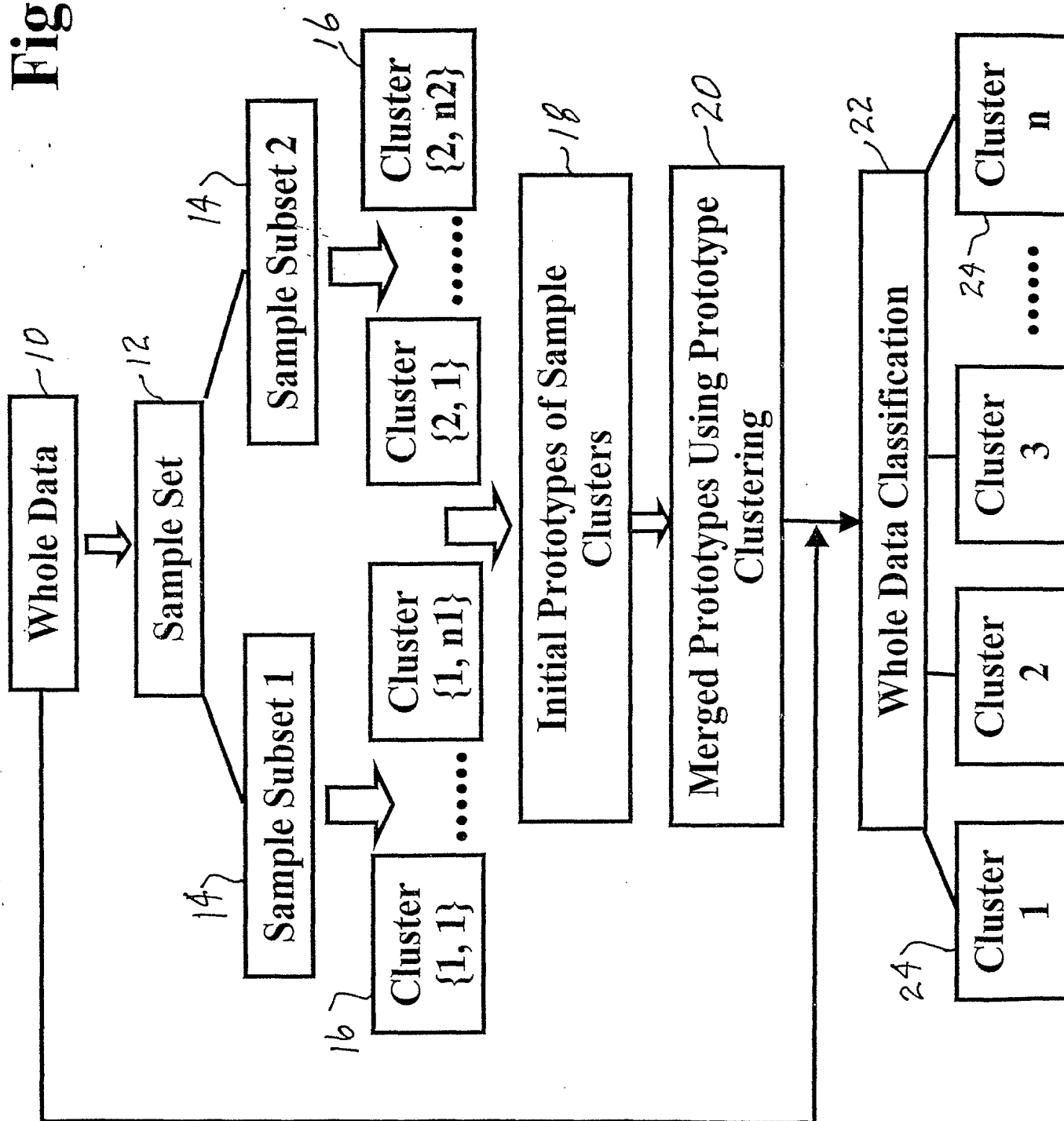


47. A method according to Claim 43, wherein the first data sample set comprises data collected by carbon oxygen logging.
48. A method according to Claim 43, further comprising:
  - using the clusters to identify geological facies types, and
  - training separate predictive models for each facies type.
49. A method of operating a hydrocarbon bearing field, comprising:
  - drilling a plurality of wells in the hydrocarbon bearing field,
  - performing open hole logging in a subset of the wells,
  - performing cased hole logging in substantially all of the wells including the subset of wells,
  - using open hole logging data and cased hole logging data from the subset of wells to train a predictive model to produce synthetic open hole data in response to inputs of cased hole data, and
  - using the trained predictive model and cased hole data from the wells to produce synthetic open hole data.
50. A method according to Claim 49, further comprising using the synthetic open hole data to plan operations for the wells.
51. A method according to Claim 49, wherein the subset of wells comprises less than one-half of the plurality of wells.
52. A method according to Claim 49, wherein the subset of wells comprises less than one-fifth of the plurality of wells.
53. Apparatus for producing synthetic values of at least one geophysical parameter for a well, comprising a predictive model trained by:
  - collecting a first data sample set of a plurality of geophysical parameters, including the at least one geophysical parameter, relating to a first well,
  - dividing the first data sample set into a set of clusters based on the geophysical parameters other than the at least one geophysical parameter,
  - selecting data from each cluster, and
  - training the predictive model to produce a synthetic value of the at least one geophysical parameter in response to inputs of the plurality of geophysical parameters, excluding the at least one geophysical parameter.

54. Apparatus according to Claim 53, wherein the predictive model comprises an artificial neural network.
55. Apparatus according to Claim 53, wherein the predictive model comprises computer code.
56. A method for producing synthetic geophysical measurements in a well having at least one depth interval in which one or more actual measurements cannot be accurately taken, comprising:
  - collecting at least one log of a plurality of geophysical measurements in a borehole, the at least one log having missing or defective measurements of at least one parameter in at least one depth interval,
  - selecting a training data set comprising at least a portion of the plurality of geophysical measurements from depth intervals other than the at least one depth interval,
  - training a predictive model with the training data set, and
  - using the trained predictive model to produce synthetic values of the missing or defective measurements of the at least one parameter in response to inputs comprising at least a portion of the geophysical measurements taken in the at least one depth interval.
57. A method according to Claim 56, wherein the at least one log comprises an open hole log and the at least one parameter is measured by the open hole log.
58. A method according to Claim 57, wherein the at least one log comprises a cased hole log.
59. A method according to Claim 58, wherein the training data set comprises parameters measured by both the open hole log and the cased hole log.

60. A program storage device readable by a machine, embodying a program of instructions executable by the machine to receive a plurality of geophysical parameters measured in a well and to produce synthetic values of at least one geophysical parameter for a well, the program comprising a predictive model trained by:
- collecting a first data sample set of a plurality of geophysical parameters, including the at least one geophysical parameter, relating to a first well,
  - dividing the first data sample set into a set of clusters based on the geophysical parameters other than the at least one geophysical parameter,
  - selecting data from each cluster, and
  - training the predictive model with the selected data to produce a synthetic value of the at least one geophysical parameter in response to inputs of the plurality of geophysical parameters, excluding the at least one geophysical parameter.
61. Apparatus according to Claim 1, wherein the predictive model comprises an artificial neural network.

Fig. 1



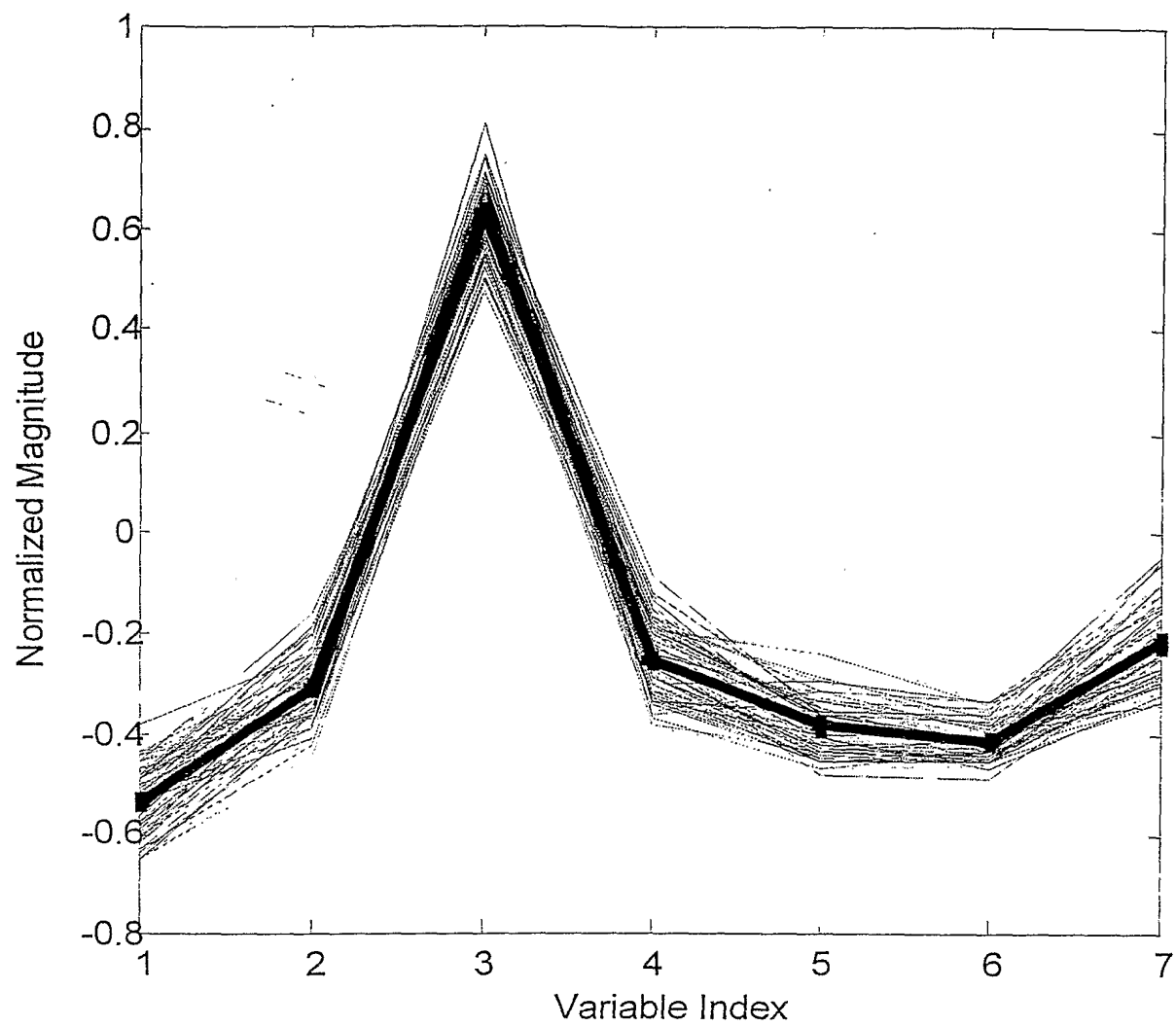


Figure 2A

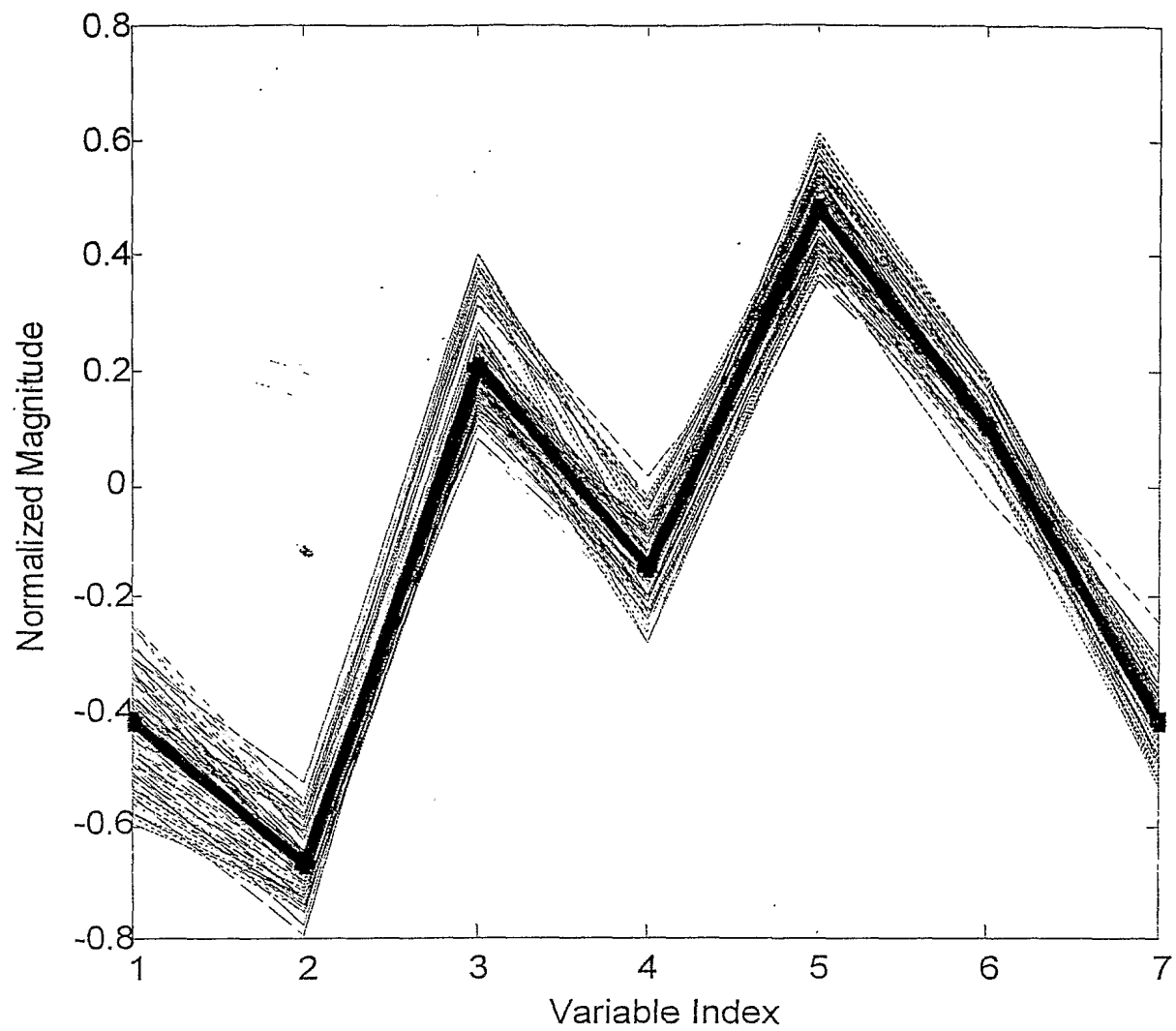


Figure 2B

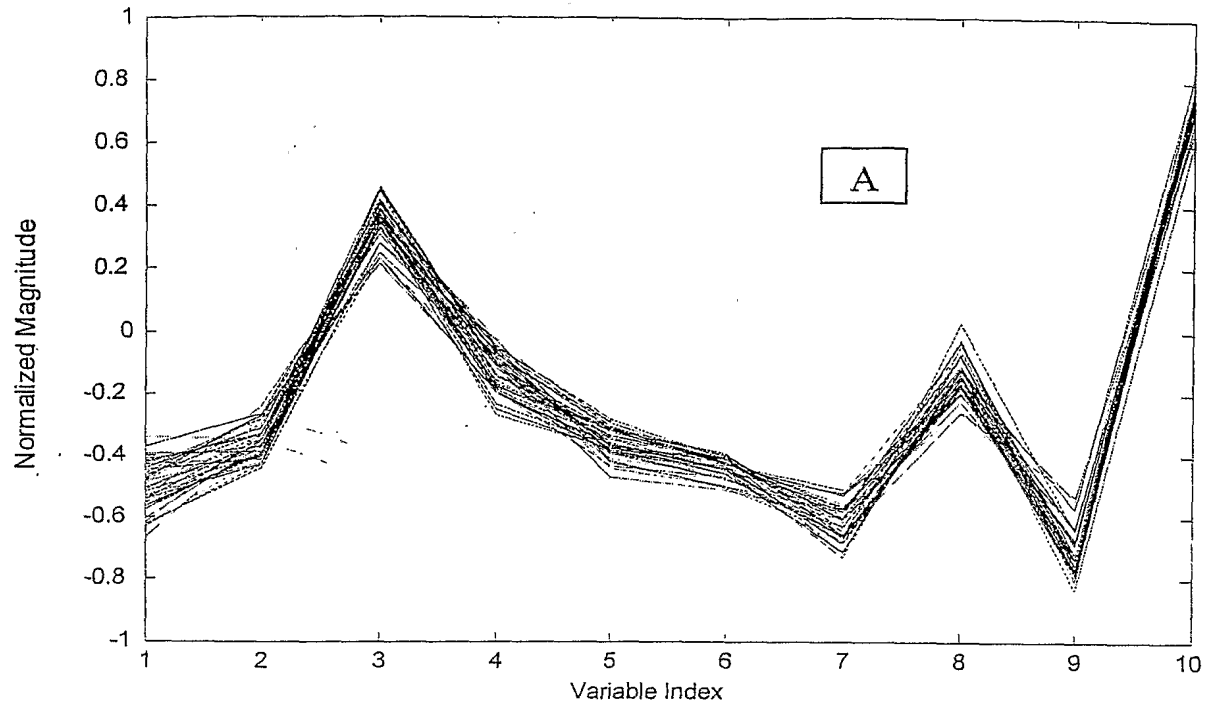


Figure 3A

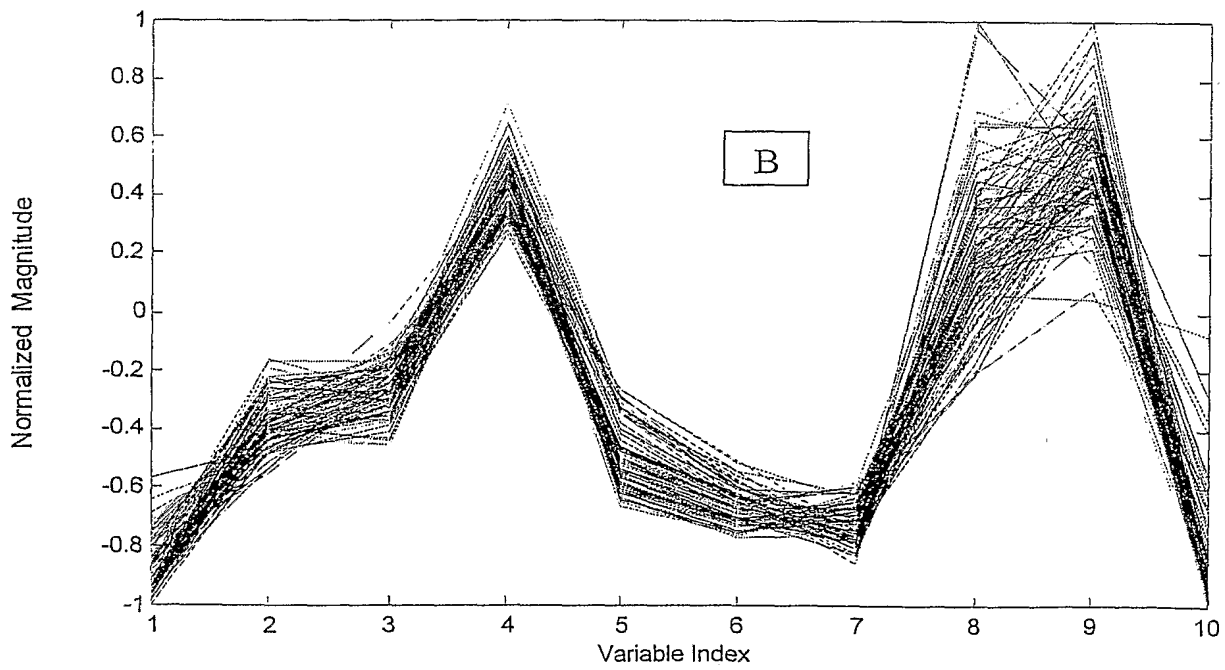


Figure 3B

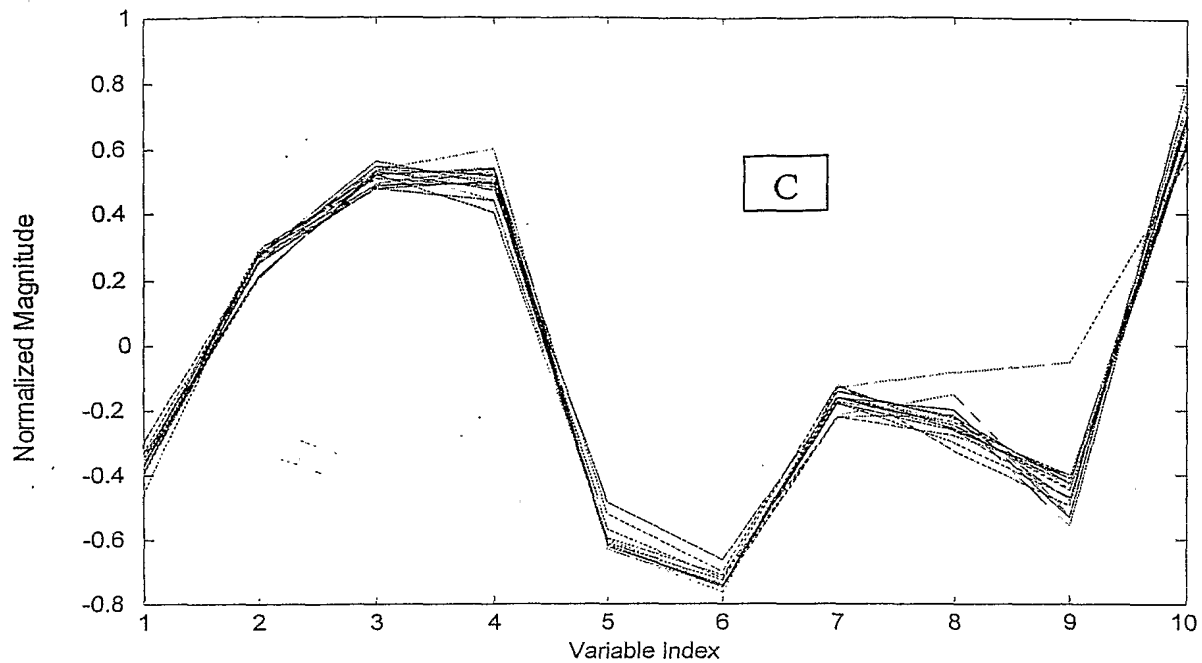


Figure 3C

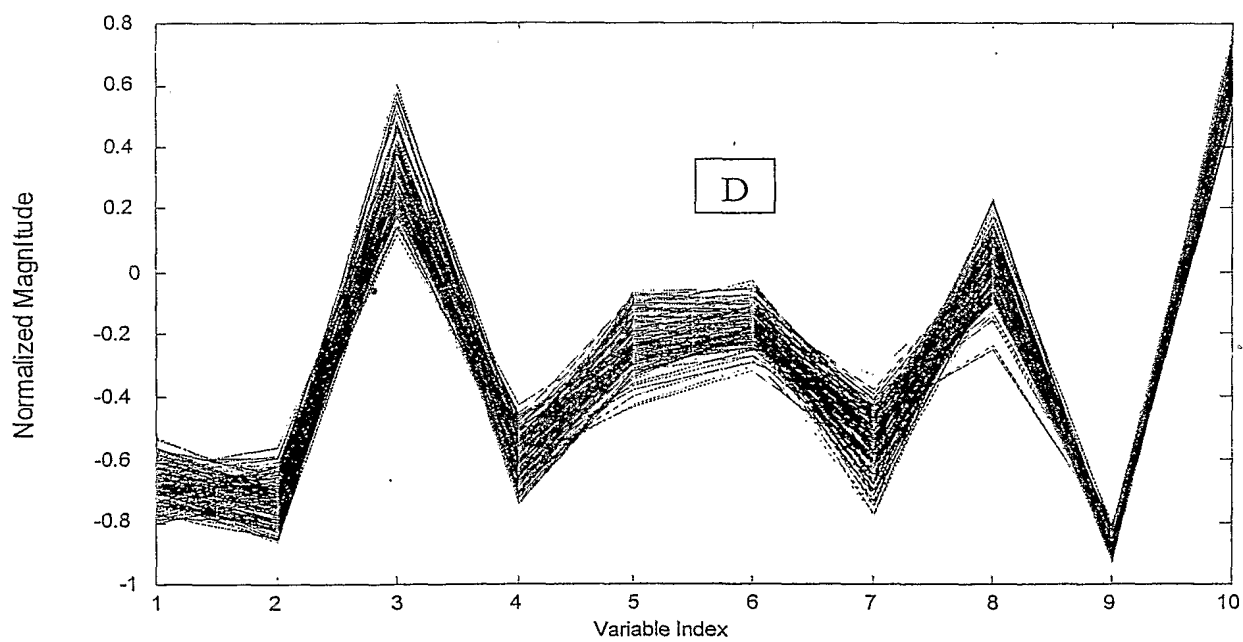
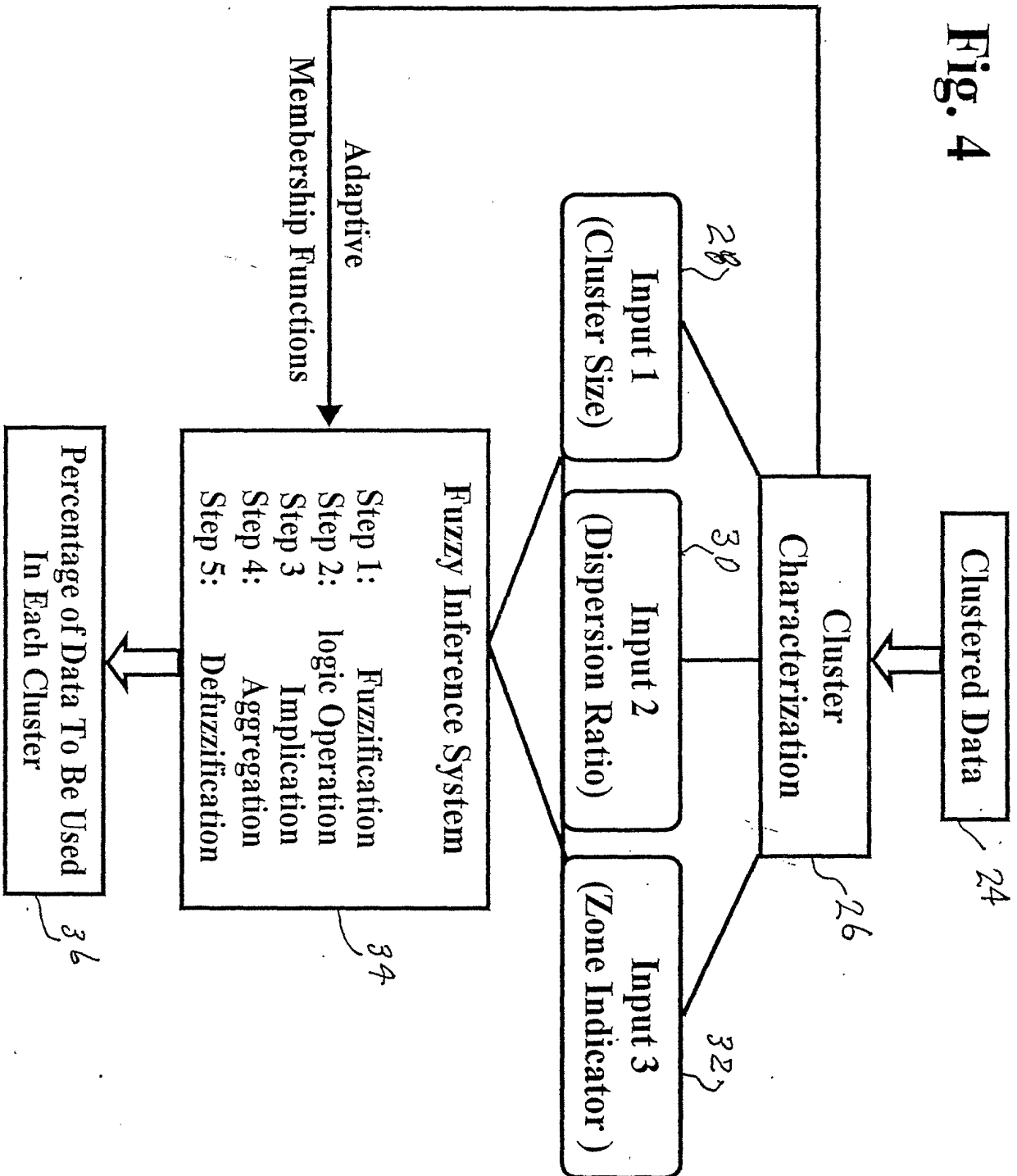


Figure 3D



Fig. 4



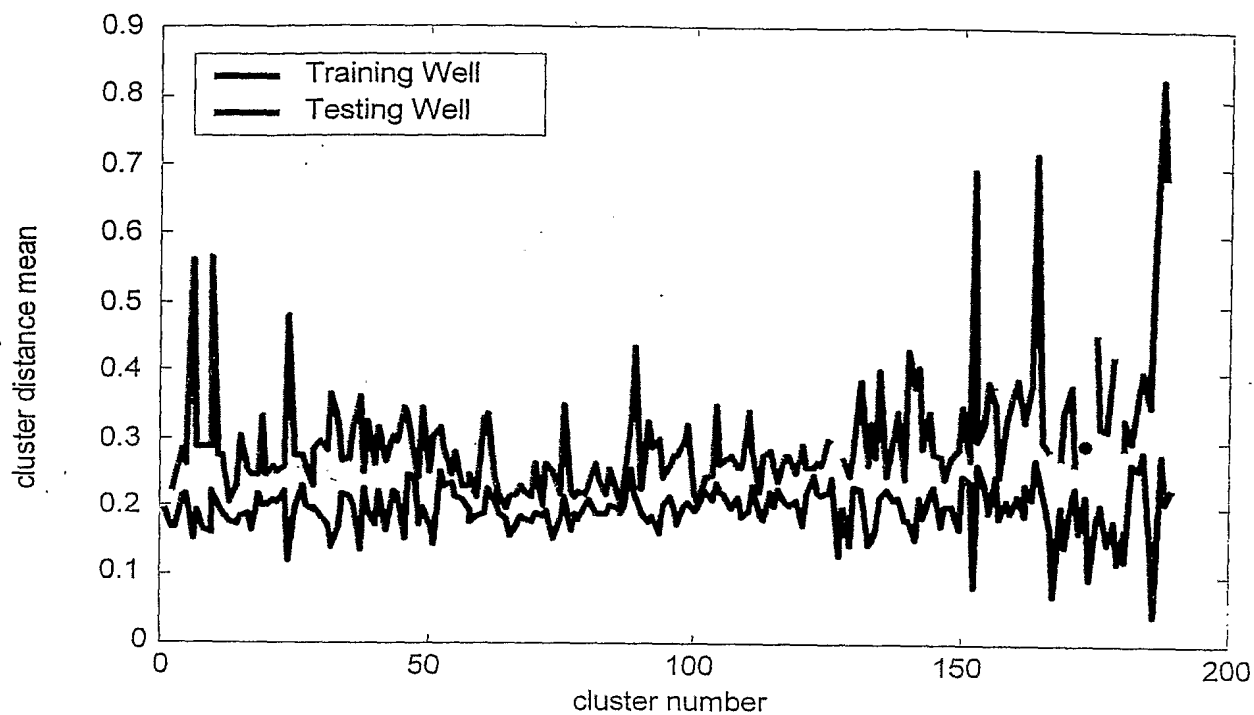


Figure 5

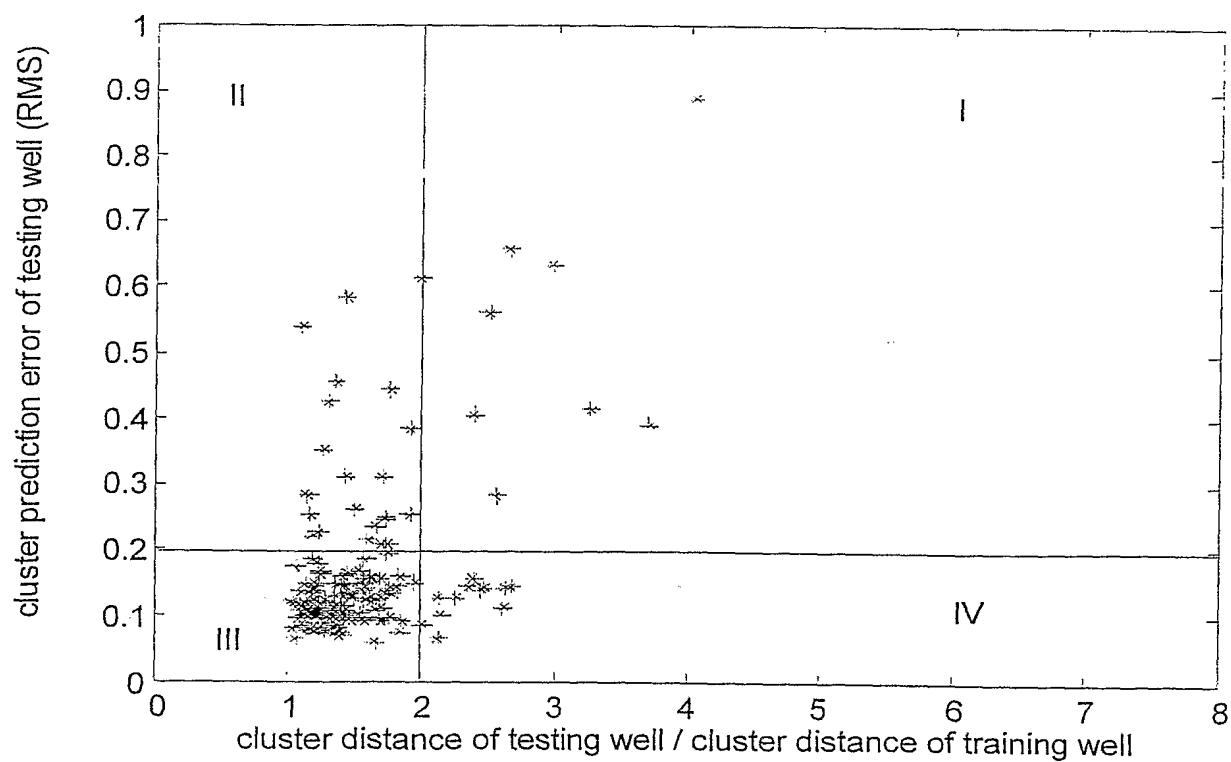


Figure 6

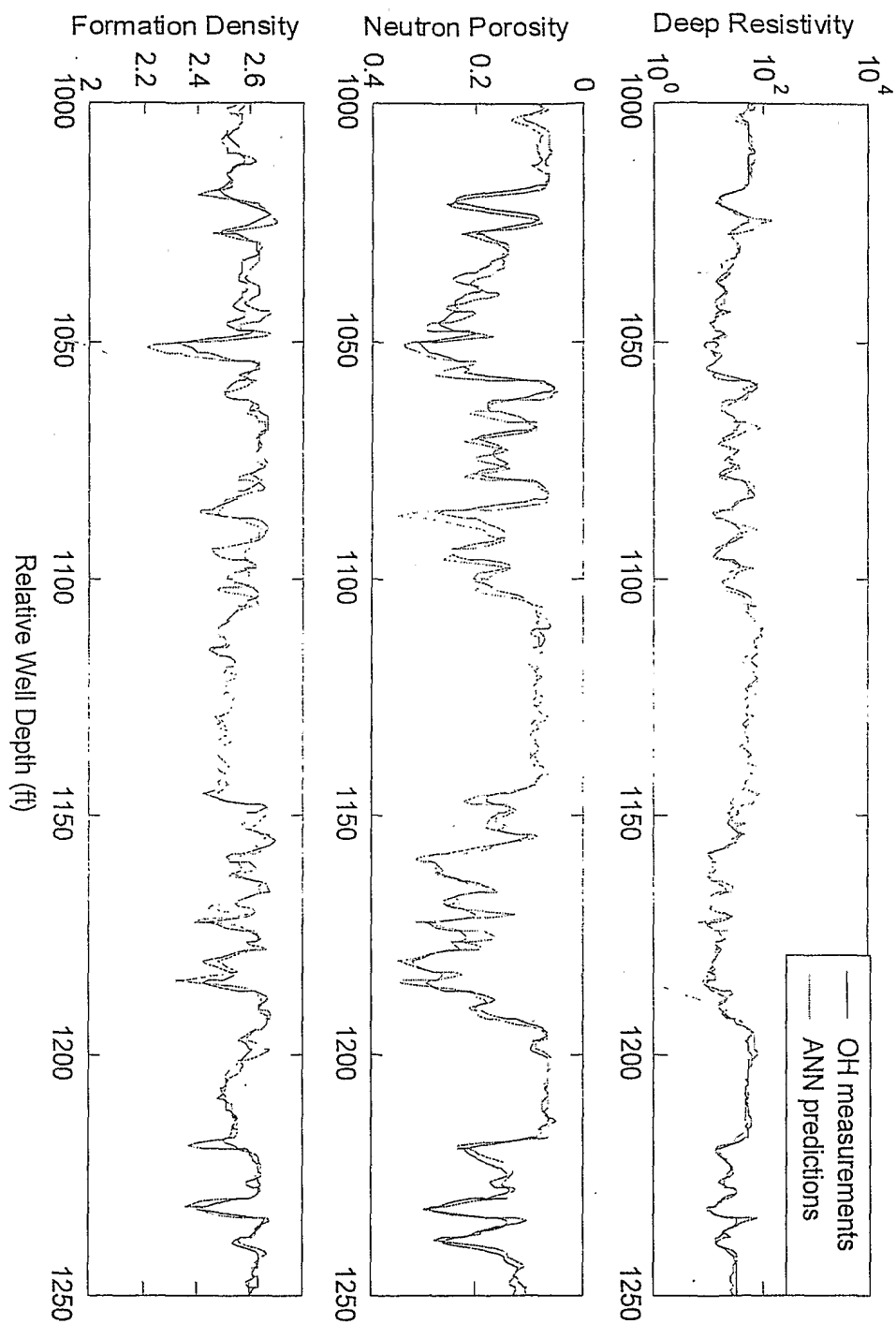


Figure 7

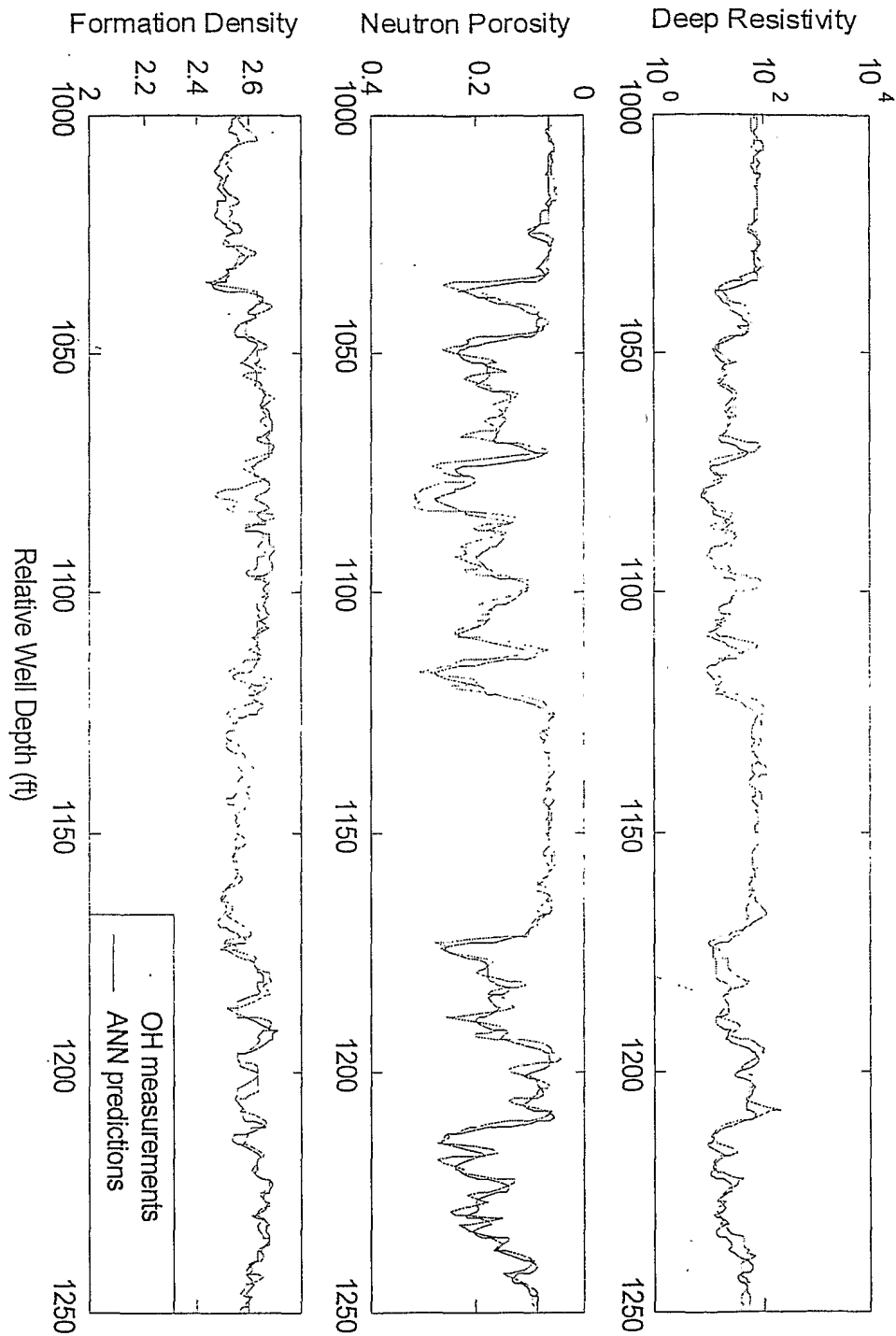


Figure 8

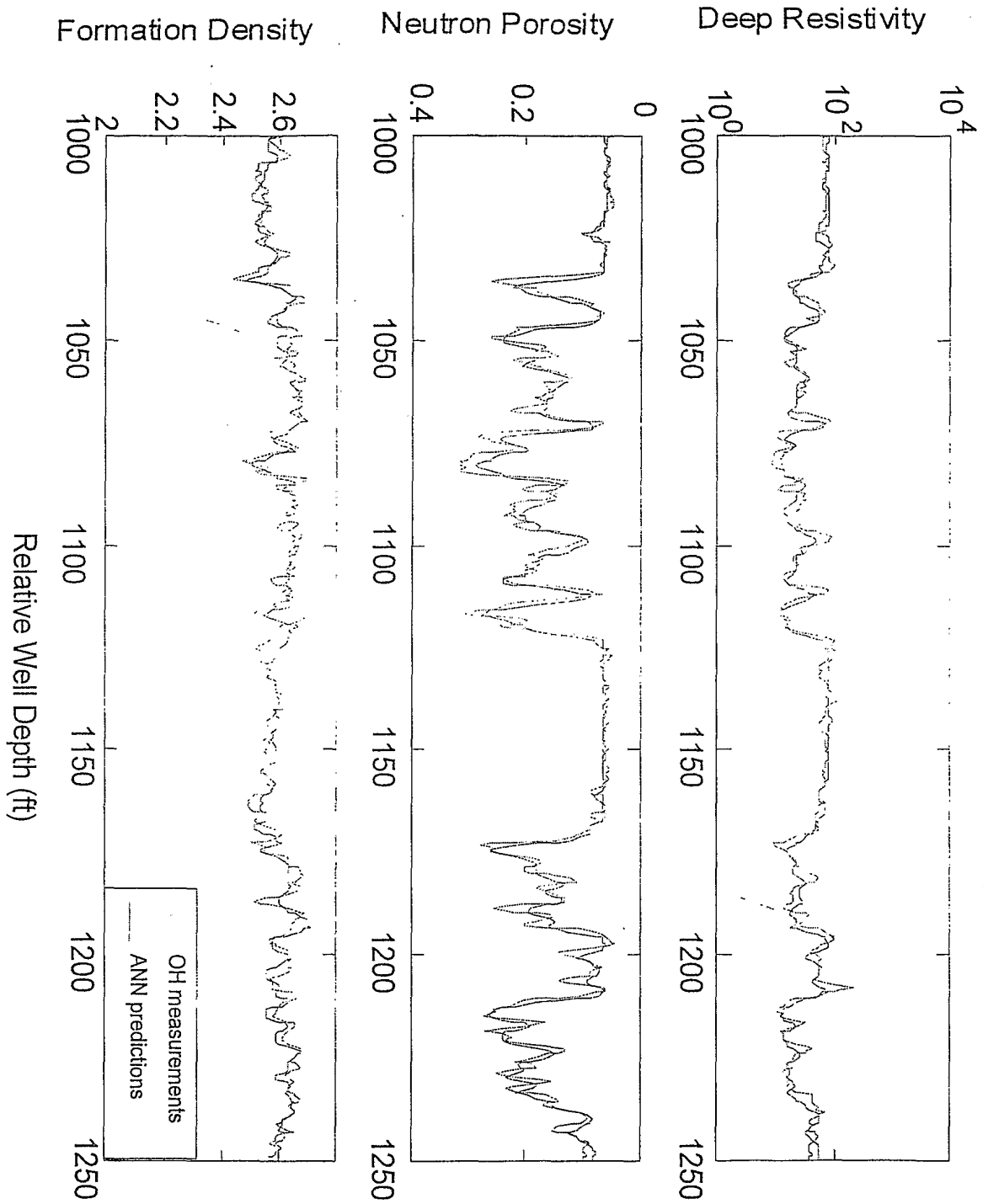


Figure 9